

Supply Chain Analytics

Vineeth Menon, MSc Big Data Analytics, L00163249

Abstract—The application of Big Data Analytics has exploded in many sectors as the amount of data available across digital platforms has increased. Like all other sectors, there have been tremendous amount of transformation in Supply Chain Management over the period due to digitization in ordering process, automated tracking of material movement, inventory etc. This has led the organization to leverage Big Data Platforms, Internet of Things and Cloud technologies to improve the efficiency in overall supply chain management. This Technical Report covers application of Big Data Analytics and Machine Learning (ML) using Databricks as a processing platform to analyze a Supply Chain dataset. The raw data set is of an E-Commerce company that deals with sporting, clothing and electronics supplies. The report explains the complete schema of datasets and explain the approach used to analyze the prime attributes. The report also explains the complete ML pipeline from selection of features essential to develop a model to determining its accuracy before deployment. All the inferences regarding the key attributes are covered in the result section that describes the accuracy of each implemented model and steps involved to improve it.

Index Terms—Apache Spark, Big Data Analytics, Light Gradient Boosting Machine, Random Forest Classifier, Decision Tree Classifier, Machine Learning

I. INTRODUCTION

With immense scope of Big Data applications in the domain of Supply Chain, this technical project explores a dataset used by DataCo Global that consists information on order, shipping and profit details for an E-Commerce company [1]. The primary motive of this project is to study the dataset and understand the rationale behind the application of data visualization and Machine Learning Classification algorithms to infer various hypothesis formulated through the dataset. The cloud service platform planned to use for executing the analytics will be Azure Databricks where as Pyspark will be the interface for Apache spark to implement Python APIs.

It is important to understand the schema of a dataset, total number of columns or dimensions involved, study distinct values in each dimensions, identifying key performance indicators and basic process knowledge before starting the processing and analysis. With 53 attributes, the project targets 'Delivery Status' of a product that can indicate if a product is delivered on time or late. It also focusses on 'Profit per Order' because the dataset covers order data, shipping details and its profit on per product in every region across the world. Hence it is essential to understand which region leads in maximum profit and minimum profit contribution. The project also focuses on studying the total 'Order Status' that is 'Pending' market wise and region wise. There are lots of other hypothesis that can be formulated but based on the definition of each attribute, these 3 have been given prime importance based on its impact in business performance improvement. The focussed goals of this project is to

1. Predict if a delivery is late or on time based on relevant features available
2. Identify the market with most pending orders. Based on it, identify the top countries in that market with maximum pending orders

3. Identify which order region and product has highest and lowest profit margin?

4. Identify the delivery status that tops the contribution in overall deliveries

For the execution of the project, it is important to understand the Big Data Platform in Apache Spark as the project is executed in Databricks using Pyspark. The project also emphasis on uses of multiple libraries and explain its impact. The attribute used for prediction is a categorical data, hence 2 classification algorithms are used. The selection of algorithms is justified through the accuracy measured post evaluation.

II. DATASET DESCRIPTION

The dataset for Supply Chain Analytics is acquired from <https://data.mendeley.com/> Mendeley Data is a secured open source cloud-based repository for storing research datasets and reports. This particular data set selected for the project was used by DataCo Global for studying the application of Big Data title in the field of smart Supply Chain. The data has various version and this project will use latest version (version 5) published on March'19. It consists of 3 csv files. DataCoSupply-ChainDataset.csv is a structured data type that consist of 53 attributes and around 180000 entries of order details covering 50 category of products segmented under sporting, clothing and electronics supplies. This file consists of vital details such as expected days for shipping, actual days for shipping, total order quantity, profit per order, shipping mode etc. Second csv file is an unstructured data type with information of product details, URL and IP address. Third csv file consists metadata information defining each attributes in above mentioned files. The definition for key dimensions in the data are as follows:

1. Days for shipping (real): Actual shipping days of the purchased product.

2. Days for shipment (scheduled): Days of scheduled delivery of the purchased product.

3. Delivery Status: Delivery status of orders: Advance shipping, Late delivery, Shipping canceled, Shipping on time.

4. Late Delivery Risk: Categorical variable that indicates if sending is late (1), it is not late (0).

5. Market: Market to where the order is delivered.

6. Order Item Quantity: Number of products per order.

7. Order Status: If the order is completed or pending or on hold.

8. Category Name: Description of the product category

9. Order Item Discount: Discount value per product category.

10. Order Profit Per Order: Profit created per order.

11. Shipping Mode: The following shipping modes are presented: Standard Class, First Class, Second Class, Same Day.

12. Order state, Region, Country, Geographic locations at which order and shipping data are monitored.

These are precise numerical and categorical dimensions in the dataset. Rest are in form of string values. All the attributes will be studied further during data pre-processing stage. Fig. 1 displays the total columns in the dataframe.

```
Out[2]: DataFrameType: string, Days for shipping (real): int, Days for shipment (scheduled): int, Benefit per order: double, Sales per customer: double, Delivery Status: string, LateDeliveryRisk: int, Category Id: int, Category Name: string, Customer City: string, Customer Country: string, Customer Email: string, Customer Name: string, Customer Id: int, Customer Name: string, Customer Password: string, Customer Segment: string, Customer State: string, Customer Street: string, Customer Zipcode: int, Department Id: int, Department Name: string, Latitude: double, Longitude: double, Market: string, Order City: string, Order Country: string, Order Customer Id: int, order date (DateOrders): string, Order Id: int, Order Item Cardinal Id: int, Order Item Discount: double, Order Item Discount Rate: double, Order Item Id: int, Order Item Product Price: double, Order Item Profit Ratio: double, Order Item Quantity: int, Sales: double, Order Item Total: double, Order Profit Per Order: double, Order Region: string, Order Status: string, Order Status: string, Order Zipcode: int, Product Card Id: int, Product Category Id: int, Product Description: string, Product Image: string, Product Name: string, Product Price: double, Product Status: int, shipping date (DateOrders): string, Shipping Mode: string
```

Fig. 1. Total columns in dataframe

III. BIG DATA IN APACHE SPARK

Apache spark is a popular open source engine. Spark provides fault tolerance as well as general purpose cluster computing system comprising Scala, Java, Python and R APIs for executing a task efficiently. Furthermore, Spark is well-suited for the creation of large scale machine learning systems because of its efficiency in iterative computational power.

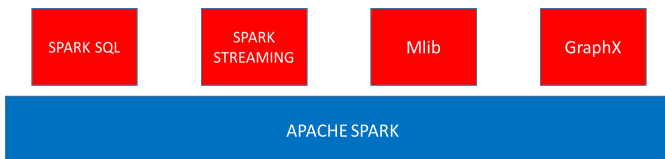


Fig. 2. Apache Spark Ecosystem

As depicted in fig2, the apache spark ecosystem consists of 4 main libraries. Spark streaming helps in real time data processing from multiple sources. GraphX is an API for graphs and graph-parallel computation. Spark SQL is a module for running SQL queries. It allows to create dataframe and act as a distributed SQL query engine. MLlib helps in implementing machine learning applications in Spark. The library aims at large-scale learning environments that can benefit from data-parallelism or model-parallelism for storing and processing data and models. MLlib provides quick and scalable applications of popular learning algorithms such as classification,

regression, collaborative filtering, clustering, and dimensionality reduction for common learning scenarios [3].

The main programming abstraction in Spark to execute parallel operations within a cluster are RDDs (Resilient Distributed Dataset). RDDs are created by applying various operations such as map, filter and groupBy functions to dataframes. RDD is a collection of elements partitioned across the nodes of the cluster that can be operated in parallel [4]. For this technical project, Databricks community version is used to create RDD. The cluster created in Databricks community edition consists of 2 cores, 15,5GB memory and 1 DBU (Databricks unit). The runtime version is 6.4 Extended Support, Spark version 2.4.5 and Scala version 2.11.

IV. CODE FILES

Git Hub link for all project documents can be found at <https://github.com/Vineeth08/BDA--Technical-Project-1.git>

V. DATA ACQUISITION AND PRE-PROCESSING

With the datasets being acquired and finalized, the primary step is to upload it to the Databricks File System (DBFS). This can be done through creating table and uploading the datasets directly from the path or from the desktop file location. This dataset has to be pulled in the form of a schema to read, write and create models. In Jupyter-Notebook, 'Pandas' library helps us to load the data in the form of dataframe. But with Apache Spark set up, the acquisition will be with the help of 'pyspark.sql.types' library. Go through below codes.

```
'from pyspark.sql.types import StructType, StructField, IntegerType, StringType, DoubleType, TimestampType
```

```
scm_df = (sqlContext.read.format("csv")'
```

StructType, StructField, IntegerType and others being imported is for Pyspark to recognize if the variable in the schema is an integer or string or float. Because of the import, the dataframe output will itself recognize the type of data in the schema. The output displays 2 jobs. The reason is that as RDD is invoked, a "job" is created. Jobs are work submitted to Spark. Jobs are divided into "stages" based on the shuffle boundary. Refer below figure. Each stage is further divided into tasks based on the number of partitions in the RDD. So tasks are the smallest units of work for Spark [7].

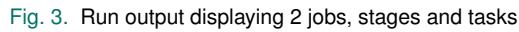
To view the schema, use command

```
display(scm_df)
```

following table with header will be displayed.

To check number of rows and columns,

```
'print((scm_df.count(), len(scm_df.columns)))'
```

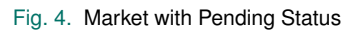


Data pre-processing is done to have a smooth executing of ML pipeline with best accuracy. Does it should not possess any null values. So it's necessary to check if current possess any null values. 3 columns possess null values. Product description – 180519 Order Zipcode – 155679 Customer Zipcode – 3 Data in all these columns are not significant to run algorithms. Hence, drop Product description, Order Zipcode. Where as using na.fill, enter '0' for 3 values of Customer Zipcode. The new dataframe is named as scm_df2.

It is important to visualize the dataframe to understand the data, its trend, pattern etc. Many questions can be answered just through mere data visualization. Good data visualization makes it easier to communicate information to larger crowd. The visualization has been given prime importance in this project. The target metrics for visualization are 'Order Status' for various markets. This will be further broken down to region wise and country wise. 'Delivery Status' region wise and 'Order Profit per Order'. All the attributes have been defined above. It is understood that with the given dataset, the most important metrics that can impact business are pending orders, late delivery and less order profit. Refer below graphs.

The fig.4 explains 'Market with Pending Status' in every market. This illustrates total orders pending as per 5 distinct market regions. Considering the organization needs to work against reducing the pending orders, a bar chart in descending order will help in prioritising the markets and its respective countries to take actions for reducing the total pending orders.

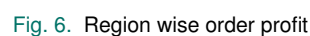
For visualization, the project imports 3 important libraries namely matplotlib, seaborn and plotly. Seaborn library is used for data visualization such as scatter plots, checking data normality etc. Matplotlib is another library used for visualization. It embeds many other visualization library (such as seaborn). It is very handy while working in multidimensional arrays. `%Matplotlib Inline` provides backend support to matplotlib. It will register a function that gets triggered whenever the output of a cell is a figure.



Country	Pending Status
Mexico	1550
Brazil	900
Honduras	450
El Salvador	400
Guatemala	380
Nicaragua	320
Panama	250
Argentina	200
Colombia	180
Venezuela	100
Peru	80
Chile	70
Paraguay	50
Ecuador	30
Bolivia	20
Uruguay	10
Costa Rica	5
Trinidad & Tobago	2
Belize	1
Suriname	1
Guatemala	1

It is learnt that Mexico and Brazil contributes to highest volume of pending orders.

Region wise order profit



Fishing category products and Cleats category products have highest profit margin, where as the organization needs to work upon increasing order profit for Soccer, Strength training and Video games going with the popularity.

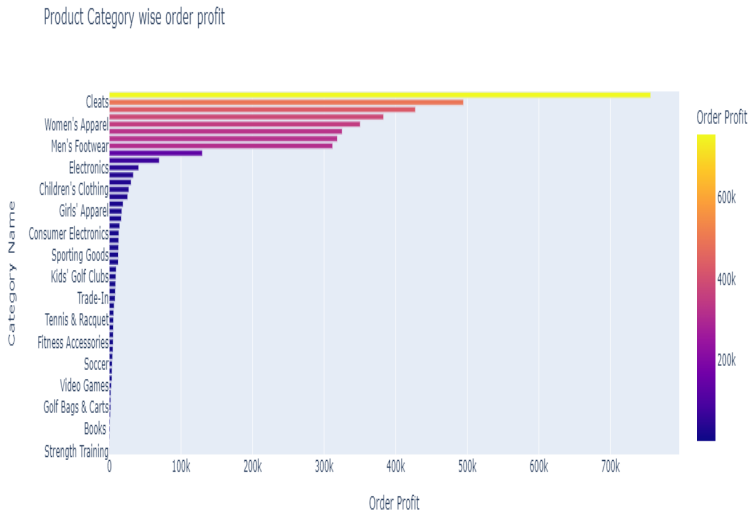


Fig. 7. Product Category wise order profit

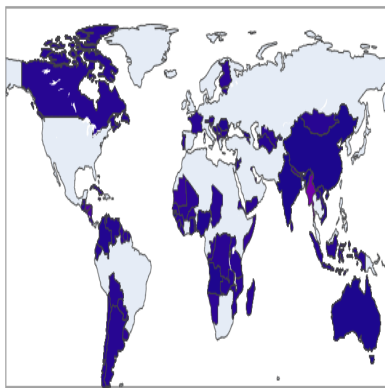


Fig. 8. Plotting order profit country wise using geographic plot

The scatter plot between Discount rate vs Sales was plotted to notice if there is any visible relation between sales to the discount rate. However, the dispersed points in scatter plot does not indicate any strong correlation.

Delivery status graph in Fig.10 is plotted to notice the impact of late delivery on overall status. It can be seen that 'Late Delivery' is a major concern for all regions and markets. Central America and Western Europe has highest late deliveries owing to the fact that number of orders are high too. Therefore, it is essential to create a model to predict the late deliveries so that it can be controlled based on vital features.

Thus through visualization, one can interpret the behaviour



Fig. 9. Scatter Plot: Order Item Discount Rate vs Sales

Delivery Status

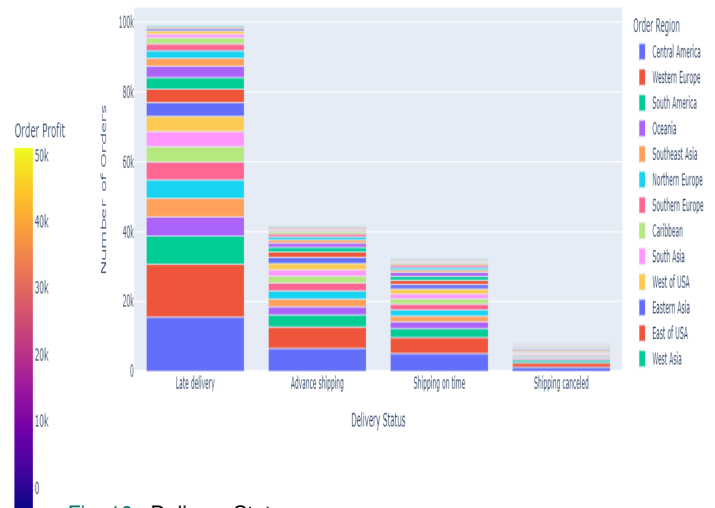


Fig. 10. Delivery Status

and patterns of key performance indicators. Post visualization and data processing, the finalized dataset can be used for developing ML models.

VII. CREATING PIPELINE AND MODEL SELECTION

The flow diagram (Fig.11) illustrates the number of steps involved in creating a complete machine learning pipeline. The data transformed earlier cannot be directly used for creating models. Identifying features to fit in a model is essential.

A. Training and Testing Dataset

Once a dataframe has been finalized after transformation, its needed to be split for training and testing. The proportion for the split is by default 80 % and 20 %, while there are different school of thoughts on selecting the actual proportion. An ideal split will avoid overfitting and under fitting of model. Some articles would tell us to use a 70:30 or even a 50:50 split. All of it depends on how large is a dataset. In this project, a trial

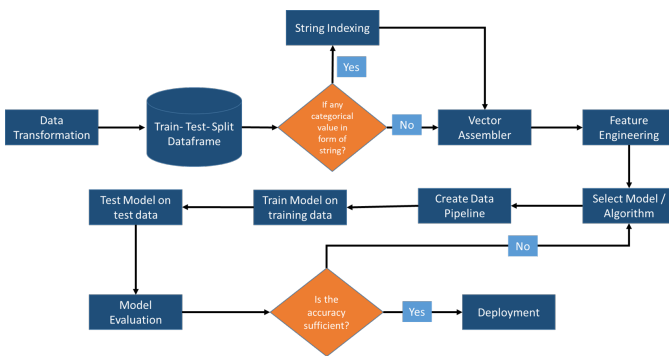


Fig. 11. Machine Learning Pipeline

and error will be taken on split ratio between 70:30 and 80:20 as per accuracy of the output.

The training data set is used to fit the model. Hence larger the data size better the model outcome. The test set is used for an fair evaluation of the final model.

B. Feature Engineering:

Feature engineering is a process where raw data is converted into features such that it can be used in machine learning algorithms. Not all data has to be transformed to features. In order to identify the exact features, domain expertise is required. The plan of model development in this dataset is to predict whether a delivery is on time or late. The label indexer used for output is the values in column `Late_Delivery_Risk`. Hence selection of features will be based on parameters that are relevant to accurate prediction of `Late_Delivery_Risk`. For feature engineering, 2 main modules needed are 'VectorAssembler' and 'StringIndexer'. Both are imported from `pyspark.ml` libraries. The role of a vector assembler is to transform multiple columns into one single vector column. The transformation is done as follows: 'class `pyspark.ml.feature.VectorAssembler(inputCols= " ", outputCol= "features")`' The name of the output column is "features". Now, some columns would have categorical data in form of strings. Machine learning algorithm will not function in presence of strings, hence it has to be converted into index or numerals using `StringIndexer` or `OneHotEncoder`. These indexed values are also added as input columns into `VectorAssembler` for feature engineering. Hence total 12 columns are added as inputs for feature column. Please note that additional features will be added and removed based on model accuracy. Refer GitHub link code files and below figure.

C. Selection of Algorithm (Model)

Selection of model depends on type of data that is needed to be predicted. In feature Engineering, feature for label column is selected as `Late_Delivery_Risk`. The supervised learning algorithms are divided into 2 categories. Regression and Classification algorithms. The main difference in application of these 2 categories are that classification algorithms are used when the value to be predicted are discrete values (labels) where as regression algorithms are used when the output is continuous data or some measurable quantity. The label

```

1 from pyspark.ml.feature import VectorAssembler
2 assembler = VectorAssembler(
3     inputCols = [
4         "Benefit per order",
5         "Sales per customer",
6         "Indexed_Type",
7         "Indexed_Market",
8         "Indexed_Shipping Mode",
9         "Indexed_Department Id",
10        "Indexed_Customer Segment",
11        "Indexed_Category Name",
12        "Order Item Discount",
13        "Order Item Discount Rate",
14        "Order Item Profit Ratio",
15        "Sales"],
16    outputCol = "features")
  
```

Command took 0.04 seconds -- by vineethmenon01@gmail.com at 2/12/2022, 4:03:28 PM on Cluster

Fig. 12. Feature Engineering using Vector Assembler

column selected for prediction has binary output. If result is 1 then delivery is late and if result is 0 then result is on time shipping or advanced shipping. The purpose of developing this model is to identify the factors (features) that will lead to late delivery shipping so that it can be controlled for least late deliveries. Since the value to be predicted is binary, this project explores possibilities of classification algorithms. In classification division, 2 algorithms selected for the project are Light Gradient Boosting Machine and Random Forest classification model. One can technically use a list of algorithms and finalise the best based on accuracy rates.

D. Light Gradient Boosting Machine

Extreme Gradient Boosting (XGBoost) and Gradient boosting decision tree (GBDT) are highly recommended classifications for analyzing the label because they offer the best mix of accuracy, stability, and computational efficiency. Light Gradient Boosting Machine (LightGBM) is an upgraded version of "gradient learning framework" based on "decision trees" and the idea of "weak" learners. Its developed by Microsoft in 2017. The difference between LGBM and XG Boost model is that the latter uses histogram-based algorithms to increase the speed up the training model, reduce memory consumption and deploys a leaf wise growth strategy with depth constraints [5].

Due to the inefficiency in terms of growth strategy, the decision tree is a poor model. It only handles the leaves of the same layer, resulting in a significant amount of wasted memory. Large number of leaf wise level can cause overfitting. LightGBM provides a "maximum depth limit" to the top of the leaf to avoid overfitting that enables high efficiency. In the project, mml Spark library is installed using Maven coordinates. MML spark library possess packages for most of the classification and regression models. Import `LightGBMClassifier`. The featuresCol will be "features" derived from `VectorAssembler` and labelCol=`Late_Delivery_Risk`. Number of iterations considered is 100.

E. Random Forest Classifier

Random forest is an example of Ensemble Learning. Ensemble learning is a method where multiple models try and solve one problem. It's an ensemble of Decision trees and decides class type through voting. Due to this Random Forest tends to

Fig. 13. Running LightGBMClassifier using Pyspark

```

Cmd 65
1 from pyspark.ml import Pipeline
2 from pyspark.ml.classification import RandomForestClassifier
3 from pyspark.ml.feature import IndexToString, StringIndexer, VectorIndexer
4 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
5
6 rf = RandomForestClassifier(labelCol="Late_delivery_risk", featuresCol="features", numTrees=10)

Command took 0.08 seconds -- by vineethnann01@gmail.com at 2/12/2022, 3:10:26 PM on Cluster

Cmd 66
1
2 rf = RandomForestClassifier(featuresCol="features", labelCol="Late_delivery_risk", numTrees=20)
3

Command took 0.04 seconds -- by vineethnann01@gmail.com at 2/12/2022, 3:10:26 PM on Cluster

```

Fig. 14. Running RandomForestClassifier using Pyspark

This is the most critical stage in the project. For evaluating the model, 'MulticlassClassificationEvaluator' has to be imported from `pyspark.ml.evaluation` library. The evaluation results for each model and further actions are covered in result section. The evaluator provides accuracy rate between the actual values and predicted values. Based on accuracy needed, correct model is selected for deployment.

Based on the prediction model run above, The LGBM Classifier provided accuracy up to 69%. This is a good accuracy but needed to be improved. Hence Random Forest Classifier model is used to recheck if the accuracy can be improved. Refer Fig. 15 and Fig. 16

The accuracy of Random Forest Classifier was also approximately 69%. Hence to check if accuracy can be improved, more features were added. 28 features are added in vector assembler that consists of 17 numerical features and 11 indexed features. The newly added features consisted of Order Regions, Product Categories, Product Names, Departments etc.

Post creating pipeline with LightGBMClassifier with all new features, the accuracy turned out to be 100%.

Even though the accuracy is best, it should be verified with other classifiers. The Random Forest needed a huge

Fig. 15. Light Gradient Boosting Machine Classifier Accuracy with ParamGridBuilder()

```
-- Job Z24 View (Stages: 1/1)

Command took 6.33 minutes -- by vimeethenon@gmail.com at 2/7/2022, 10:08:26 AM on Cluster1

Out 44

1 pred1 = model.transform(testdf)

+ ▶ pred1: spark.sql.dataframe.DataFrame = [Type: wing, Days for shipping (real): integer, ... 59 more fields]

Command took 3.52 seconds -- by vimeethenon@gmail.com at 2/7/2022, 11:07:47 AM on Cluster1

Out 45

1 from pyspark.sql.evaluation import MulticlassClassificationEvaluator
2 evaluator = MulticlassClassificationEvaluator(labelCol="Late_delivery_risk", predictionCol="prediction")
3 accuracy = evaluator.evaluate(pred1)
4 print("Accuracy = %g" % accuracy)

+ () Spark Jobs

Accuracy = 0.688943

Command took 2.43 minutes -- by vimeethenon@gmail.com at 2/7/2022, 11:08:24 AM on Cluster1

Out 46
```

Fig. 16. Random Forest Classifier Accuracy

```

1 from pyspark.ml.feature import VectorAssembler
2 feature_list = ['Indexed_Type',
3 'Days for shipping (real)',
4 'Days for shipment (scheduled)',
5 'Benefit per order',
6 'Sales per customer',
7 'Indexed_Category Name',
8 'Indexed_Customer City',
9 'Indexed_Customer Country',
10 'Indexed_Customer Segment',
11 'Department Id',
12 'Latitude',
13 'Longitude',
14 'Indexed_Market',
15 'Order Item Discount',
16 'Order Item Discount Rate',
17 'Order Item Product Price',
18 'Order Item Profit Ratio',
19 'Order Item Quantity',
20 'Sales',
21 'Order Item Total',
22 'Order Profit Per Order',
23 'Indexed_Order Region',
24 'Indexed_Order State',
25 'Indexed_Order Status',
26 'Indexed_Product Name',
27 'Product Price',
28 'Product Status',
29 'Indexed_Shipping Mode']
30
31 assembler = VectorAssembler(inputCols=feature_list, outputCol="features")

```

Fig. 17. Feature Engineering with additional attributes

```

1 from pyspark.ml import Pipeline
2 pipeline = Pipeline(stages=[IndexerInput, marketIndexer, shippingModelIndexer, late_delivery_riskIndexer, CategoryNameIndexer, CustomerSegmentIndexer,
3                             Customer_City_riskIndexer, Customer_Country_riskIndexer, Order_Ration_riskIndexer, Order_Status_riskIndexer, Product_Name_riskIndexer,
4                             assembler, classifer])
5 model = pipeline.fit(training)

+ (3) Spark Jobs

Command took 37.43 seconds -- by vrvachanmishra@gmail.com at 2/12/2022, 3:10:28 PM on Cluster

Out 75

1 pred2 = model.transform(testset)

+ (4) pred2 -> pyspark.sql.dataframe.DataFrame = (Type: string, Days for shipping (real) (neg) ... 65 more fields)

Command took 0.44 seconds -- by vrvachanmishra@gmail.com at 2/12/2022, 3:10:28 PM on Cluster

Out 76

1 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
2 evaluator = MulticlassClassificationEvaluator(labelCol='late_delivery_risk', predictionCol='prediction')
3 accuracy = evaluator.evaluate(pred2)
4 print('Accuracy = %g' % accuracy)

+ (5) Spark Jobs

Accuracy = 1

Command took 39.83 seconds -- by vrvachanmishra@gmail.com at 2/12/2022, 3:10:28 PM on Cluster

```

Fig. 18. Light Gradient Boosting Machine Classifier Accuracy with new features

